

# Parallel Performance Visualization on the Cray XT4

**Luiz DeRose**  
**Programming Environment Director**  
**Cray Inc.**  
**[ldr@cray.com](mailto:ldr@cray.com)**

**NERSC**  
**September 18-20, 2007**



Luiz DeRose ([ldr@cray.com](mailto:ldr@cray.com)) © Cray Inc.

# Cray Apprentice<sup>2</sup>

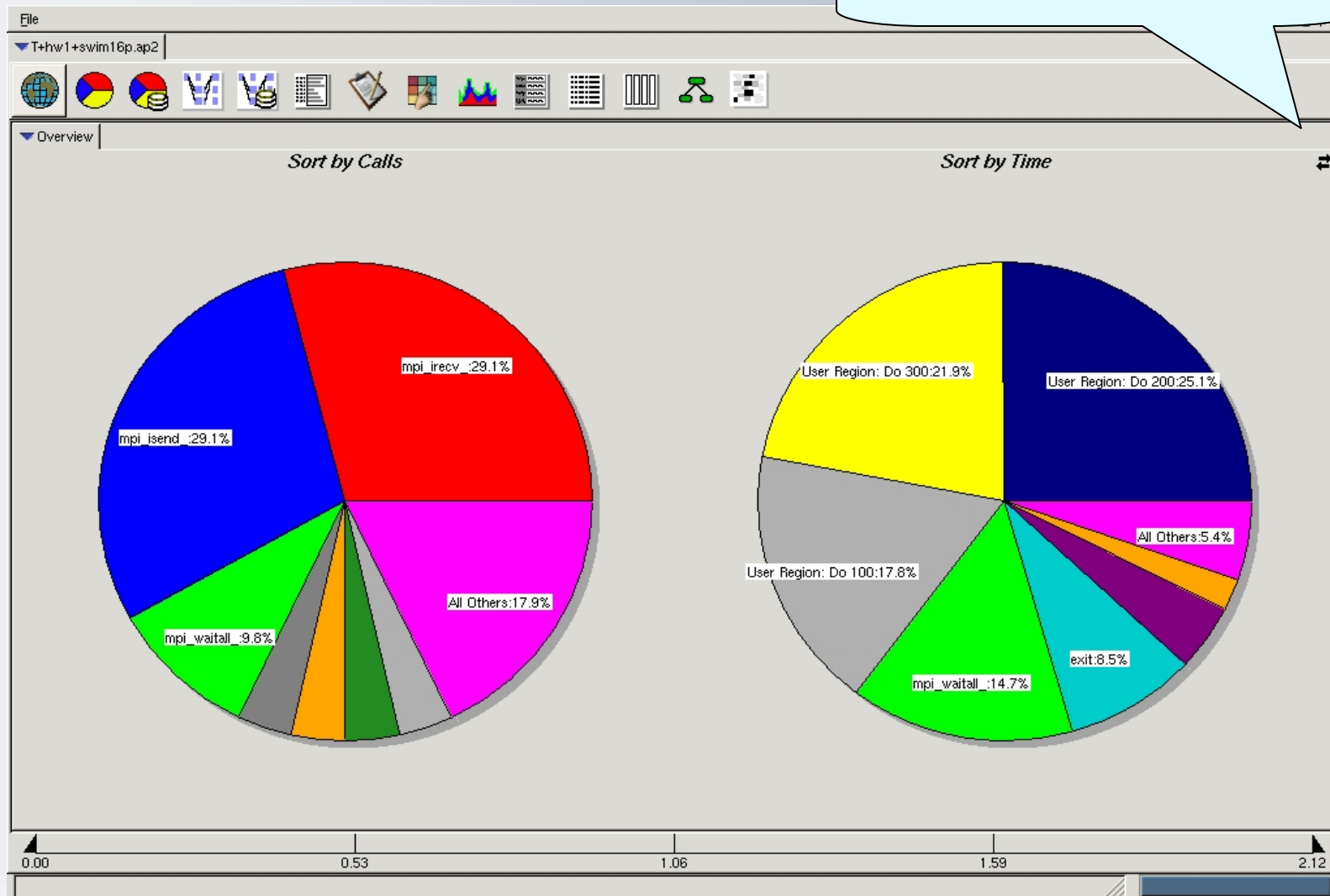
- Call graph profile
- Communication statistics
- Time-line view
  - Communication
  - I/O
- Activity view
- Pair-wise communication statistics
- Text reports
- Source code mapping

- Cray Apprentice<sup>2</sup>
- is target to help and correct:
  - Load imbalance
  - Excessive communication
  - Network contention
  - Excessive serialization
  - I/O Problems

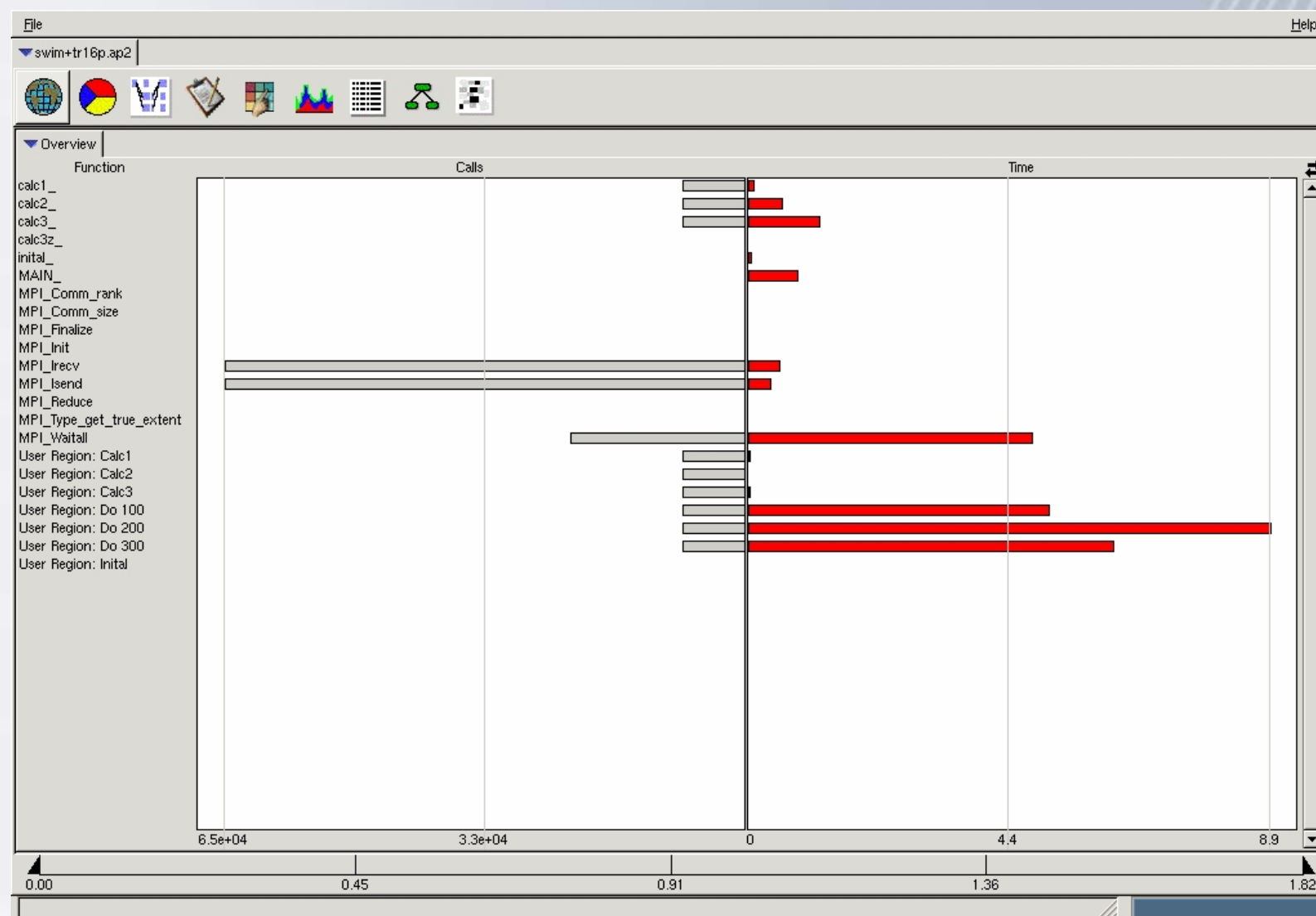


# Statistics Overview

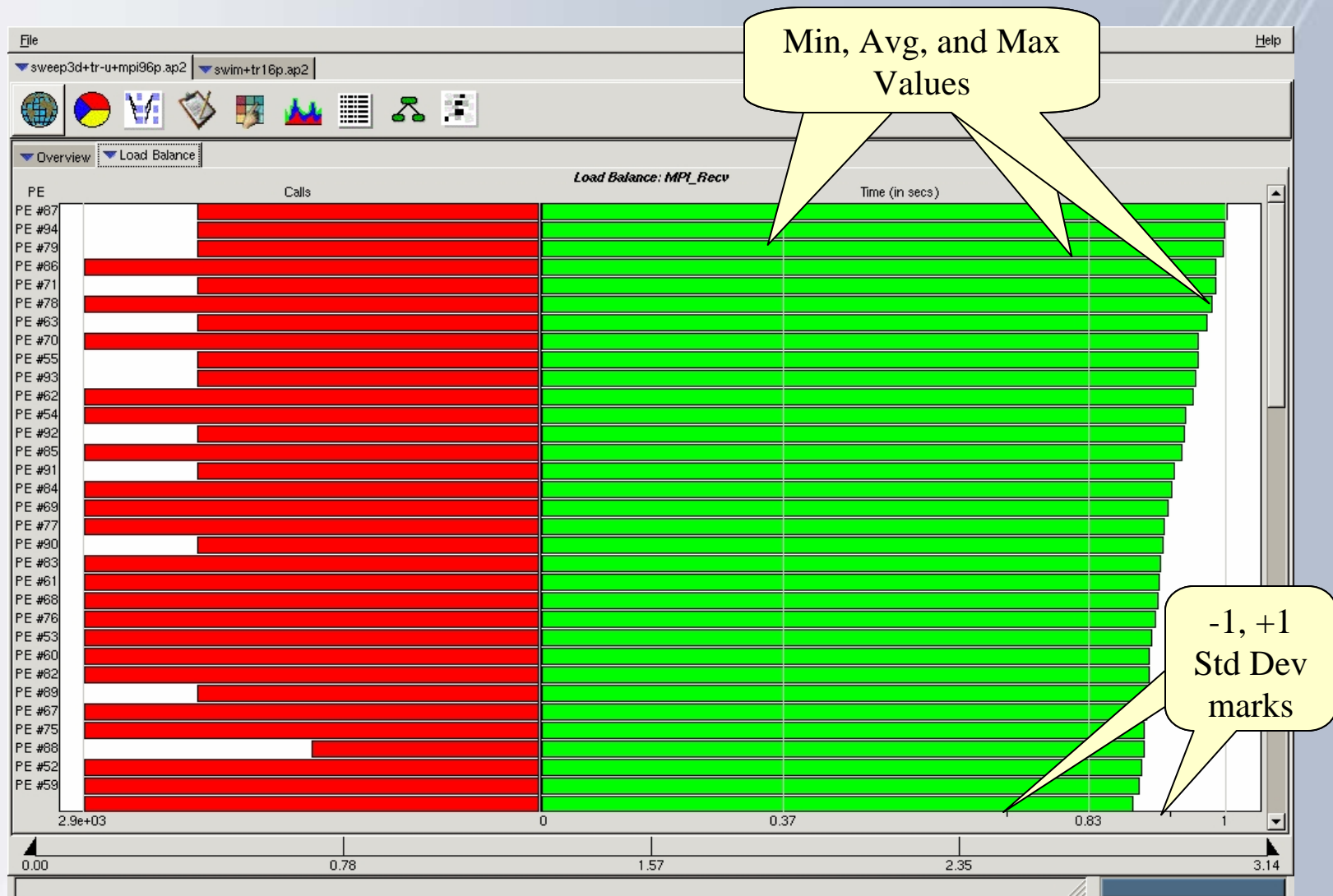
Switch Overview display



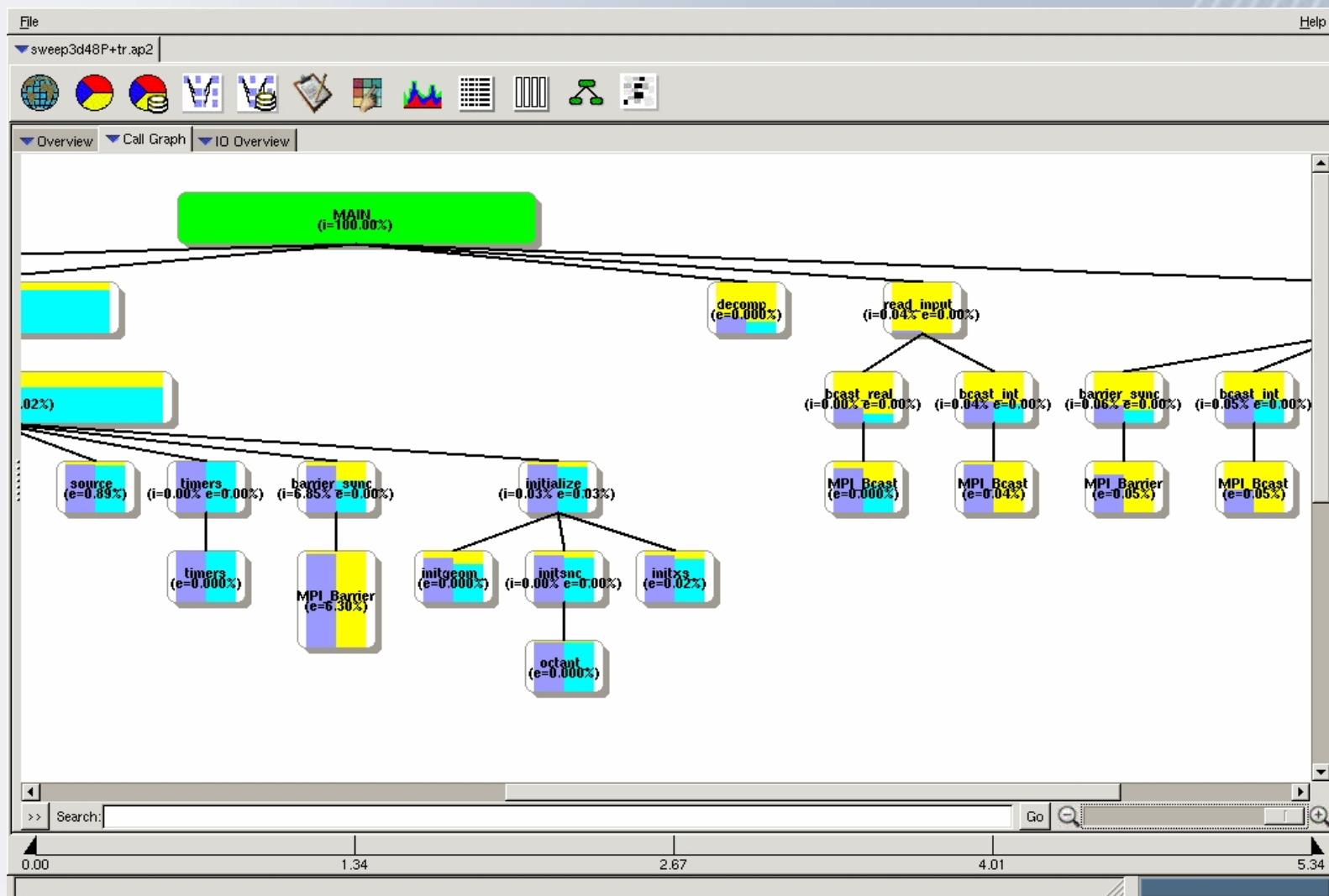
# Function Profile



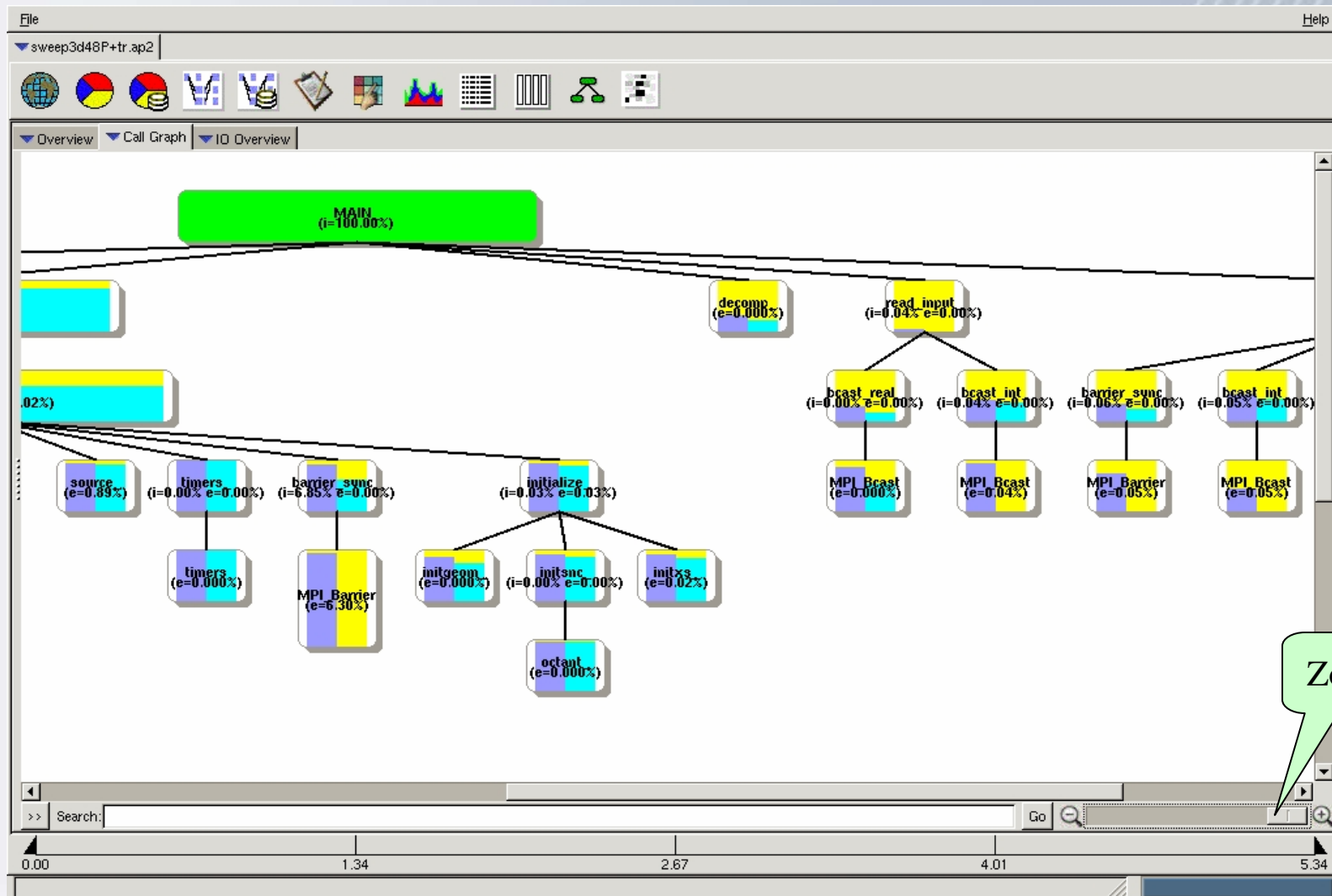
# Load Balance View (Aggregated)

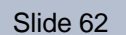


# Call Graph View



# Call Graph View

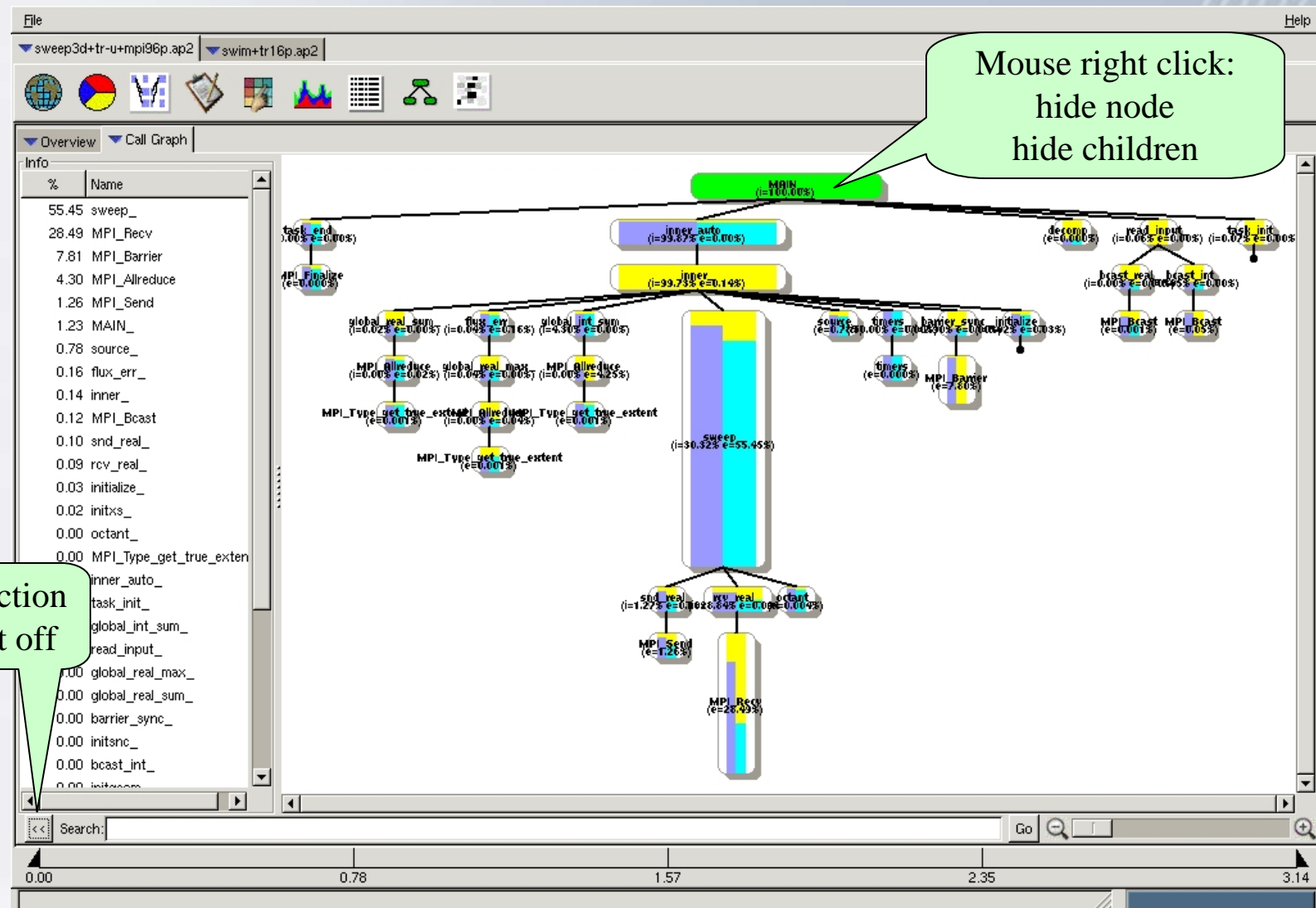




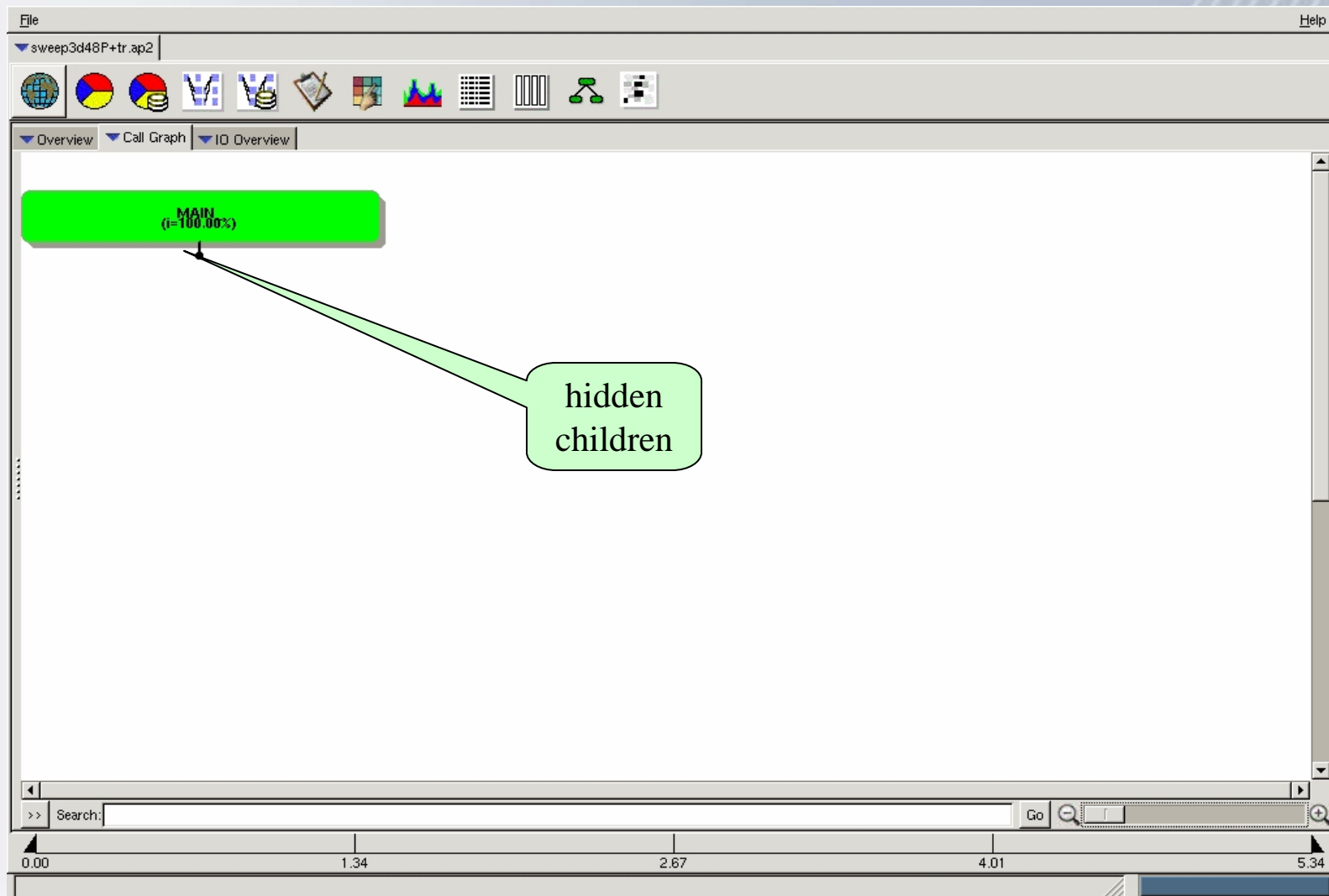




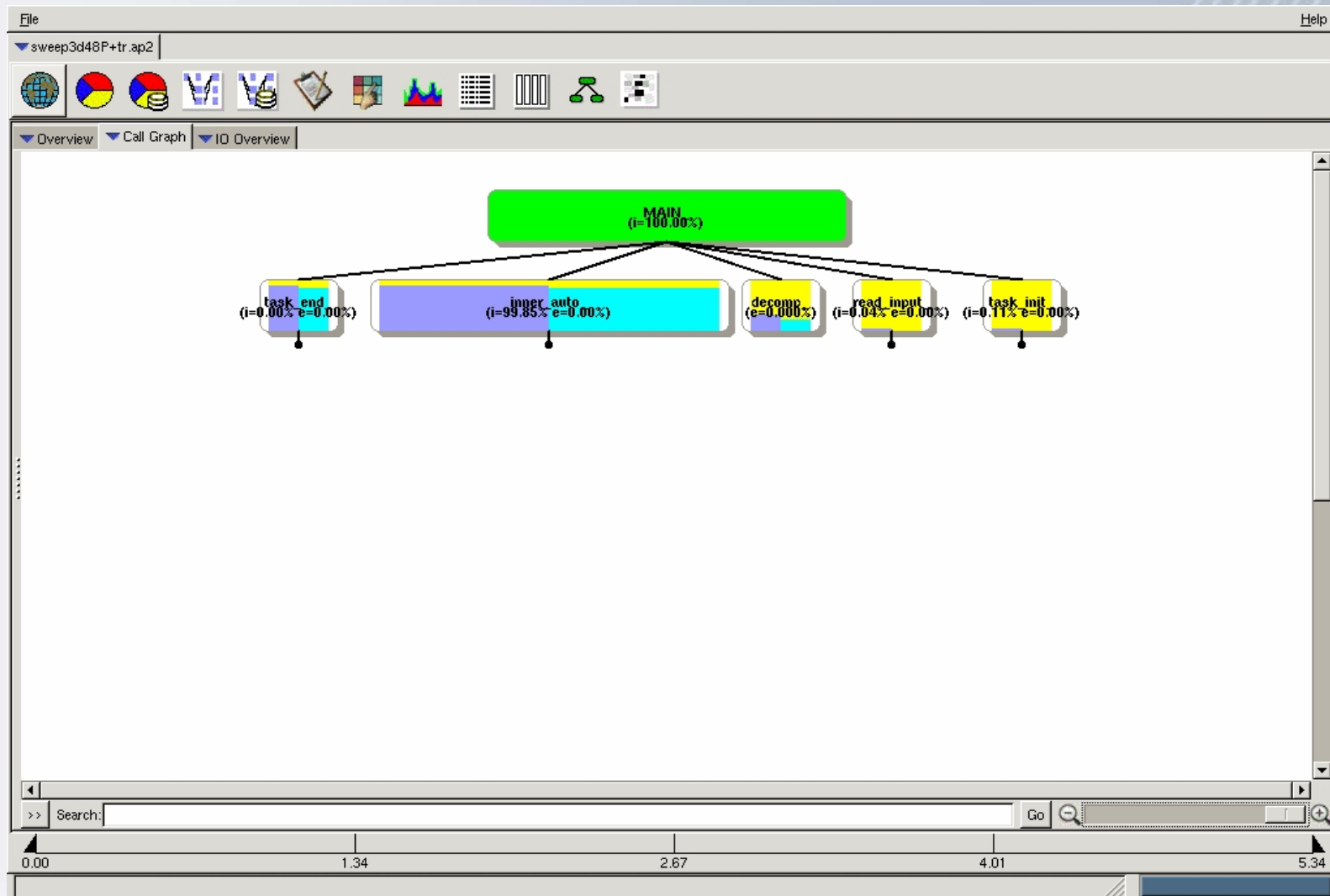
# Call Graph View – Function List



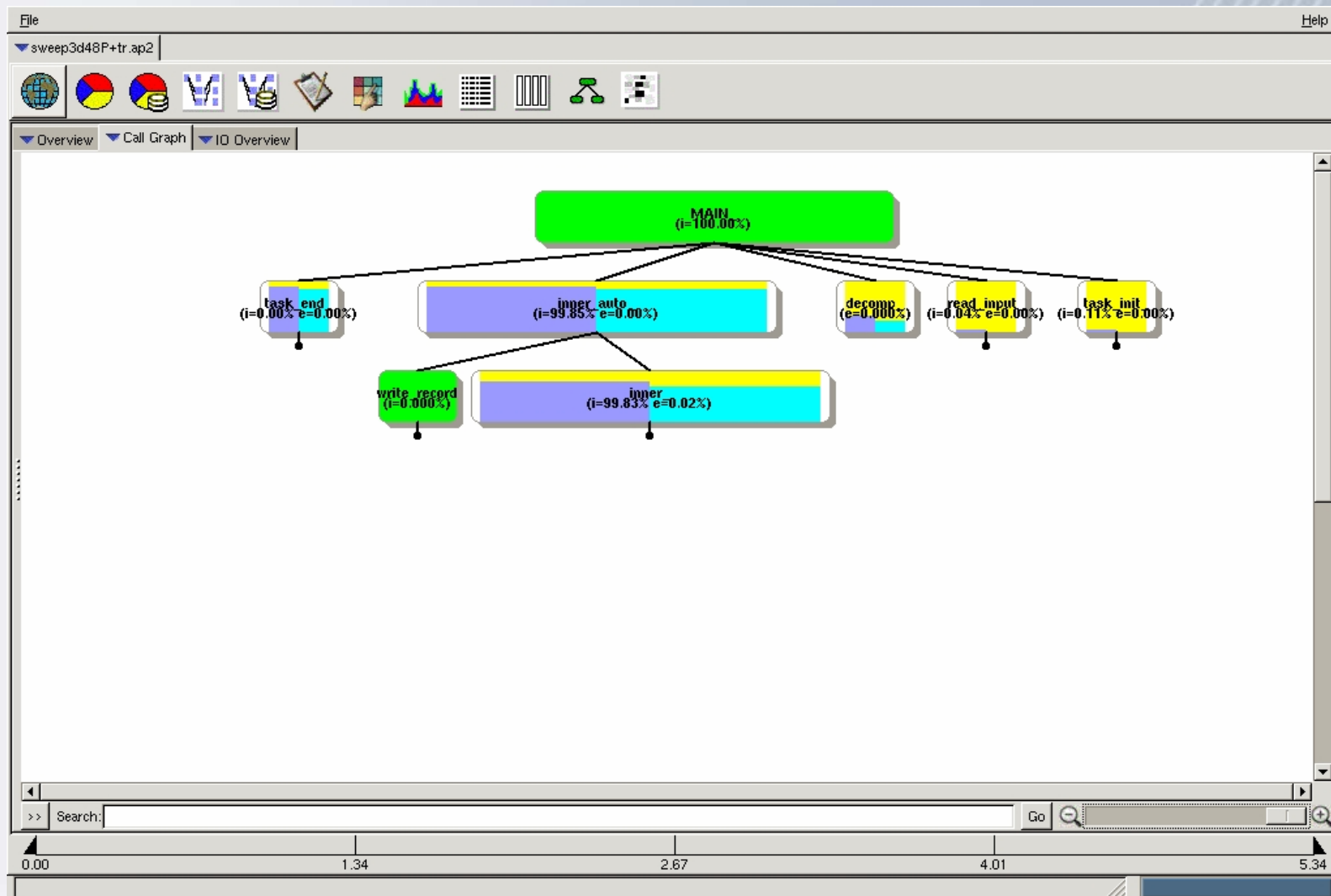
# Call Graph Hide Children



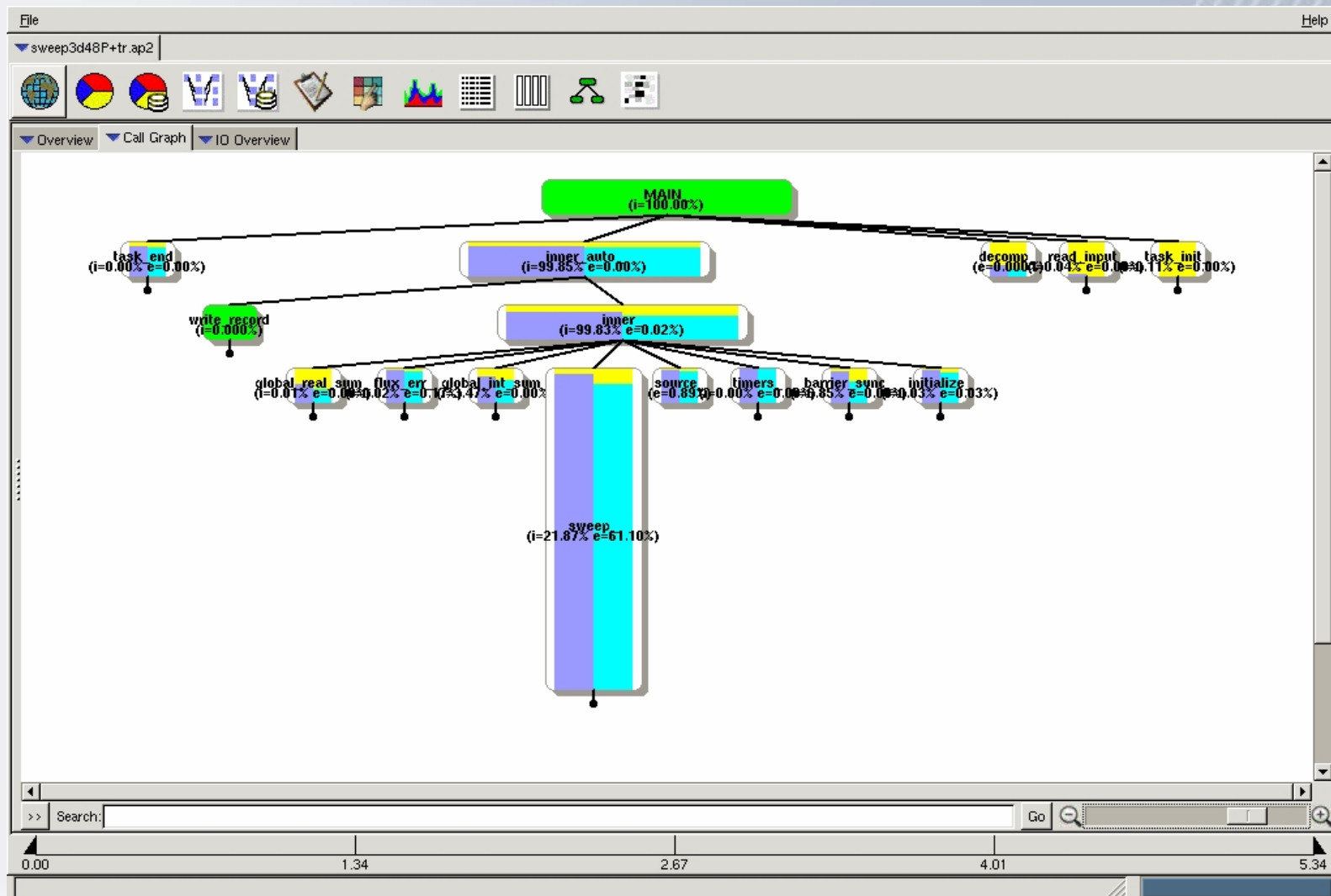
# Call Graph Unhide One Level



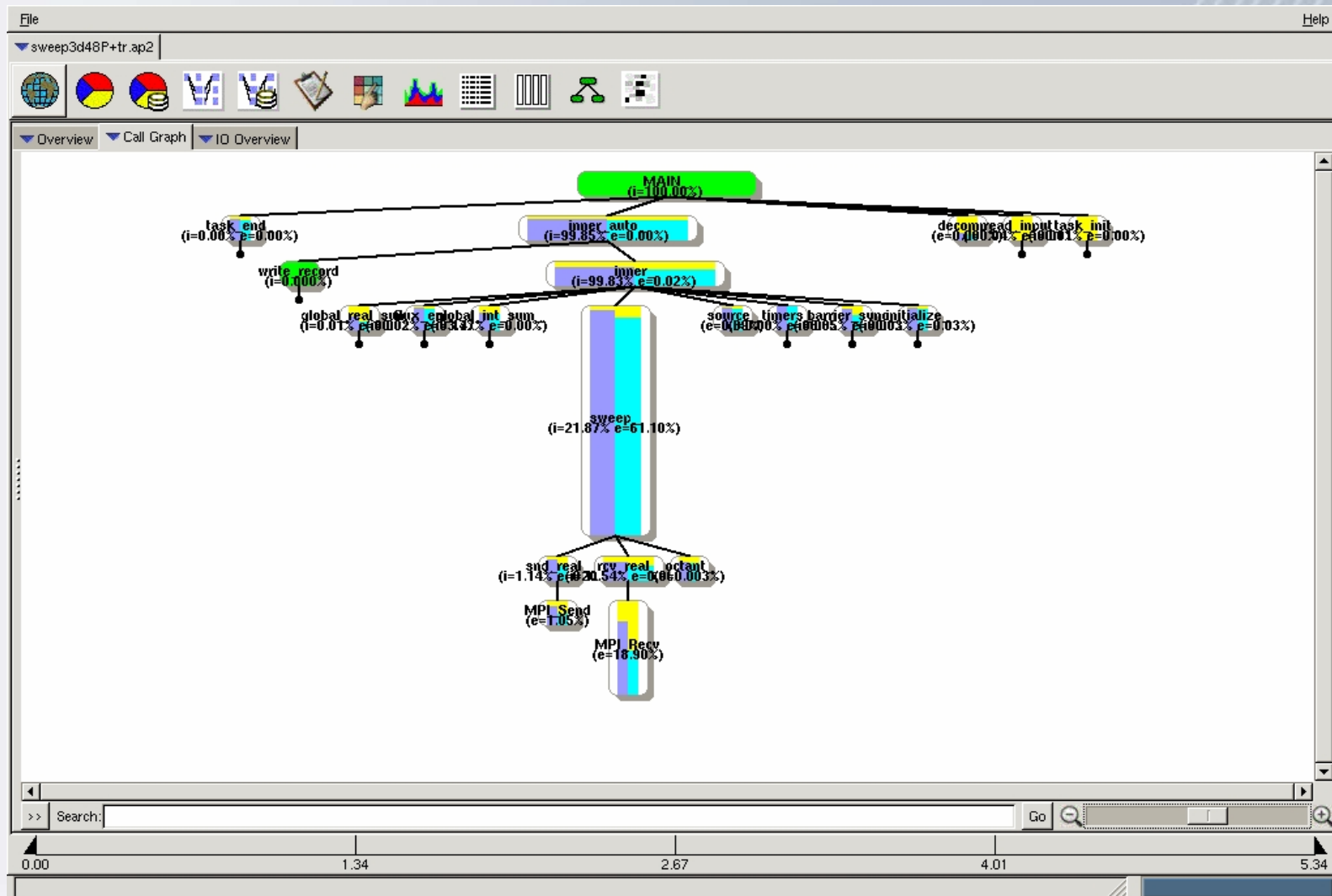
# Call Graph Unhide One Level (2)



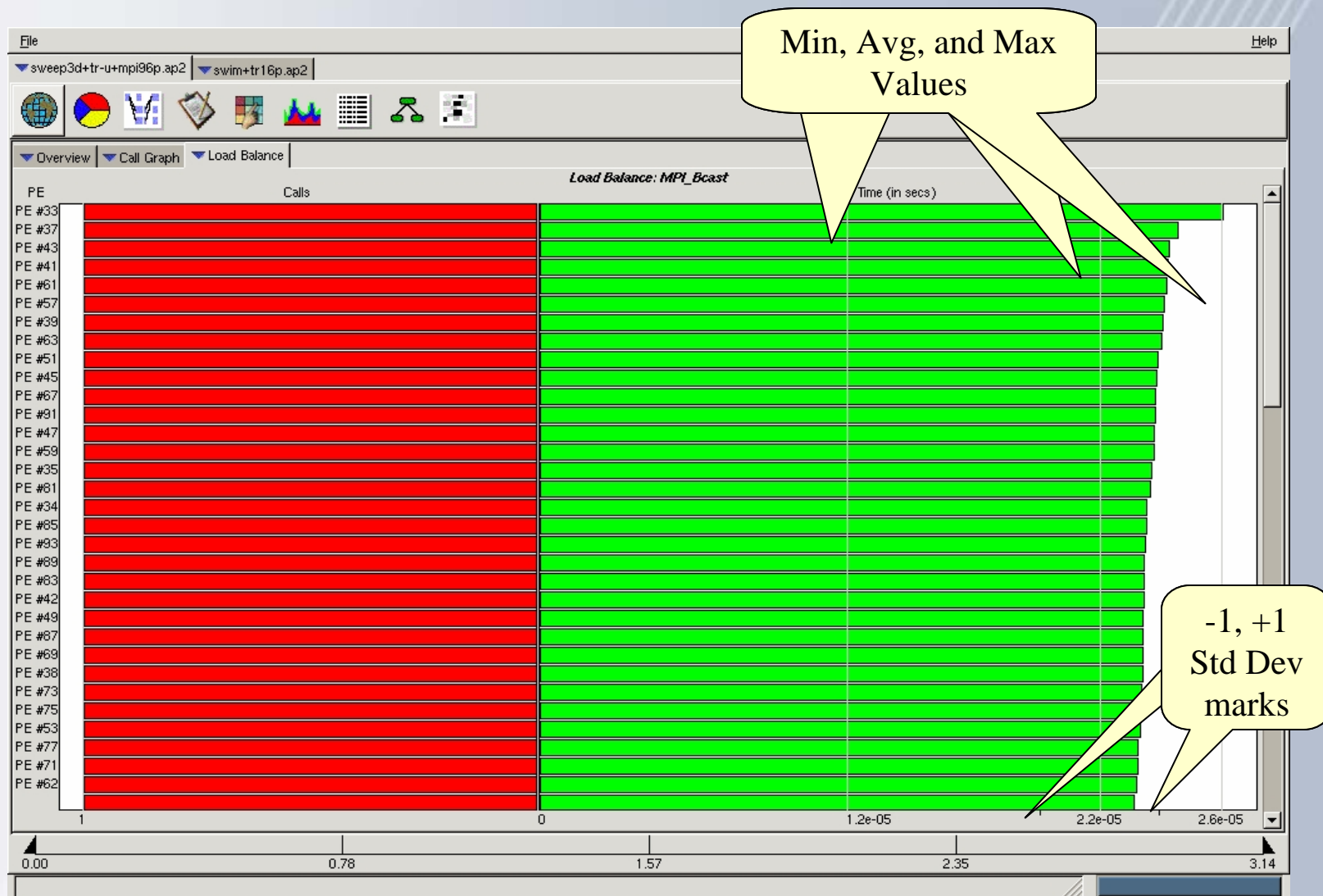
# Call Graph Unhide One Level (3)



# Call Graph Unhide All Children

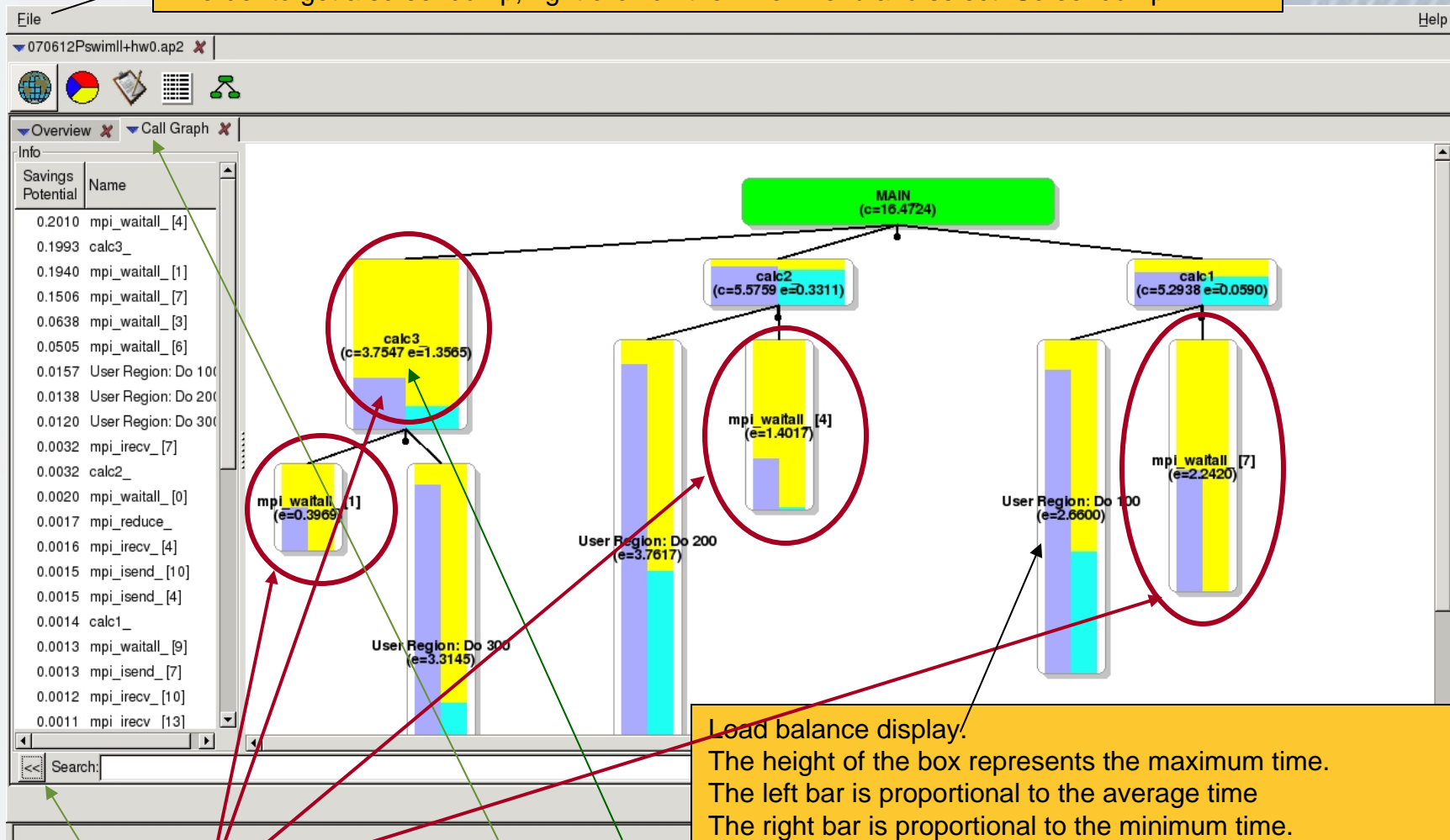


# Load Balance View (from Call Graph)





In order to get a screendump, right click on the "File" menu and select "Screendump"

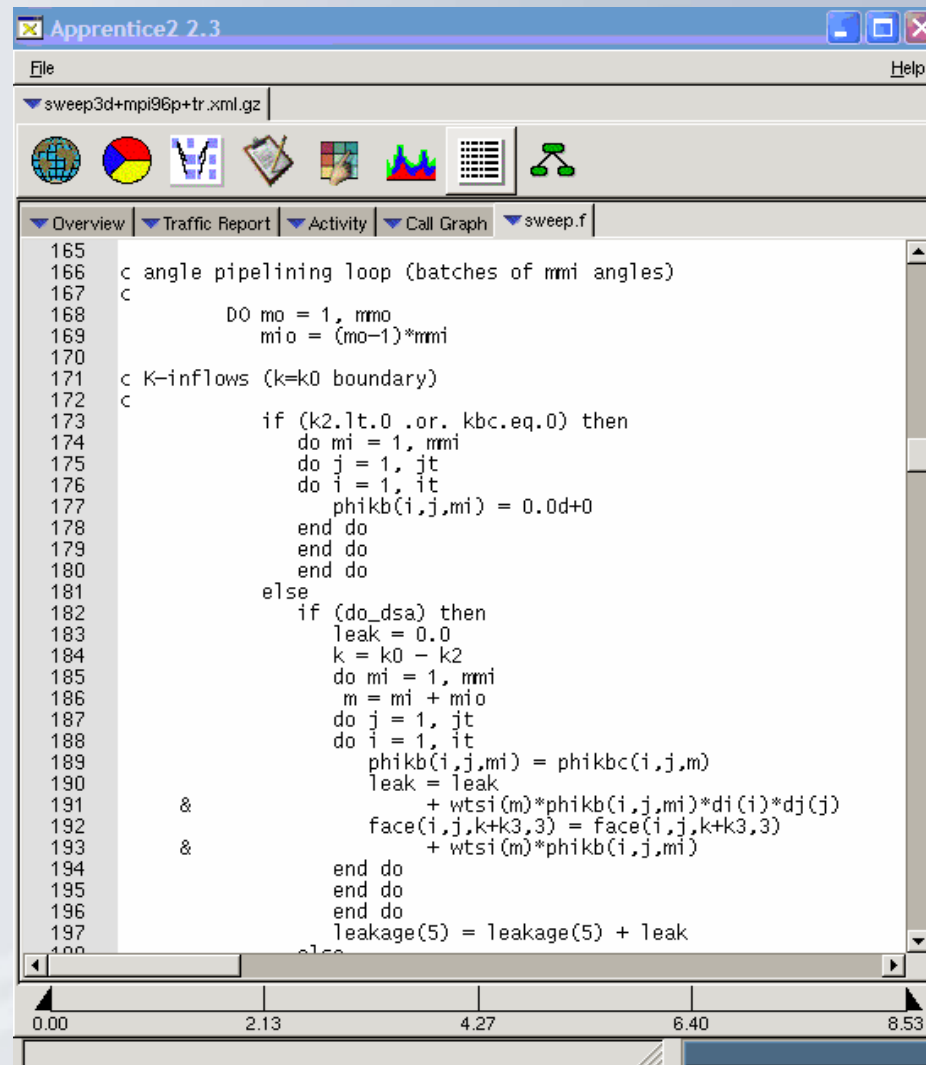


Load balance display.  
The height of the box represents the maximum time.  
The left bar is proportional to the average time  
The right bar is proportional to the minimum time.  
Functions with load imbalance have a lot of yellow on the left side

Main functions showing load imbalance. Normally load balance in MPI functions is caused by some load imbalance in computation. In this case, on "Calc3", which has an artificial imbalance.

List of functions and regions sorted by load imbalance (potential savings). In order to get this list in the call graph, click on the ">>" button and right click the call graph tab and select "potential savings".  
This call graph is filtered to show only the functions that take a reasonable amount of time.

# Source Mapping from Call Graph

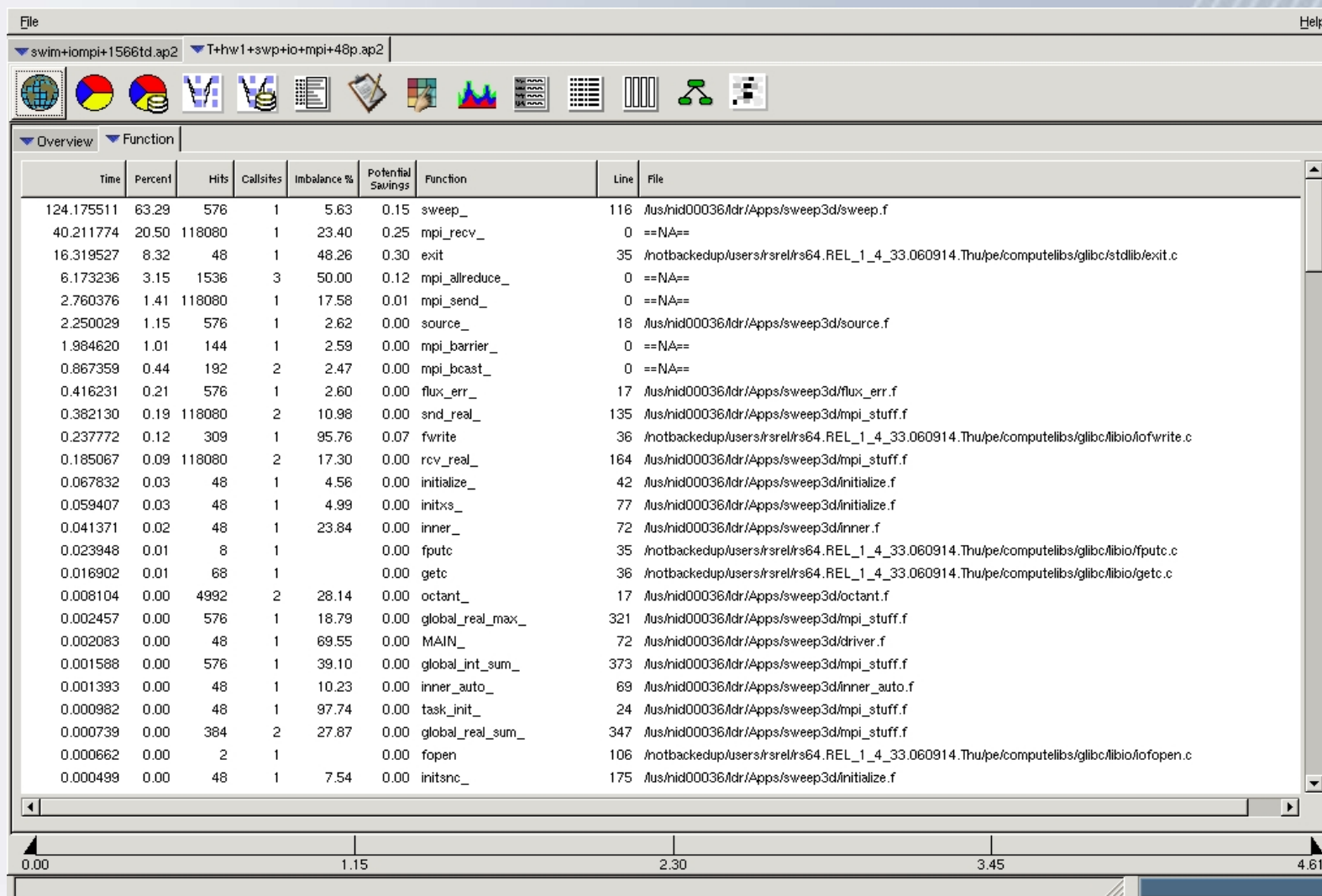


```

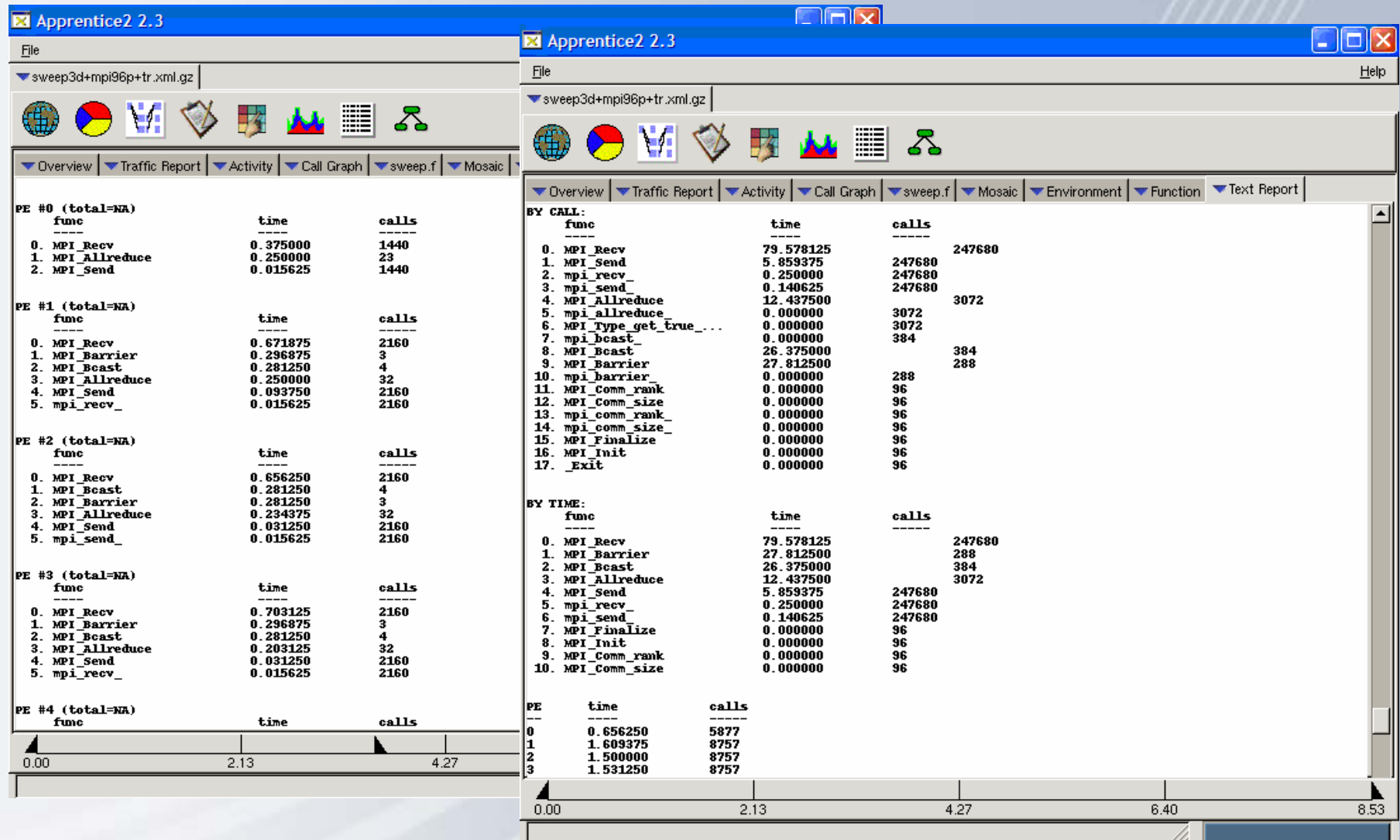
165
166 c angle pipelining loop (batches of rmi angles)
167 c
168     DO mo = 1, rmo
169         mio = (mo-1)*rmi
170
171 c K-inflows (k=k0 boundary)
172 c
173     if (k2.lt.0 .or. kbc.eq.0) then
174         do mi = 1, rmi
175             do j = 1, jt
176                 do i = 1, it
177                     phikb(i,j,mi) = 0.0d+0
178                 end do
179             end do
180         end do
181     else
182         if (do_dsa) then
183             leak = 0.0
184             k = k0 - k2
185             do mi = 1, rmi
186                 m = mi + mio
187                 do j = 1, jt
188                     do i = 1, it
189                         phikb(i,j,mi) = phikbc(i,j,m)
190                         leak = leak
191                             & + wtsi(m)*phikb(i,j,mi)*di(i)*dj(j)
192                         face(i,j,k+k3,3) = face(i,j,k+k3,3)
193                             & + wtsi(m)*phikb(i,j,mi)
194                     end do
195                 end do
196             end do
197             leakage(5) = leakage(5) + leak
198         else
199

```

# Function Profile



# Distribution by PE, by Call, & by Time



# Environment & Execution Details

File Help

swim+impi+1566td.ap2 T+hw1+swp+io+mpi+48p.ap2

Overview Function Environment

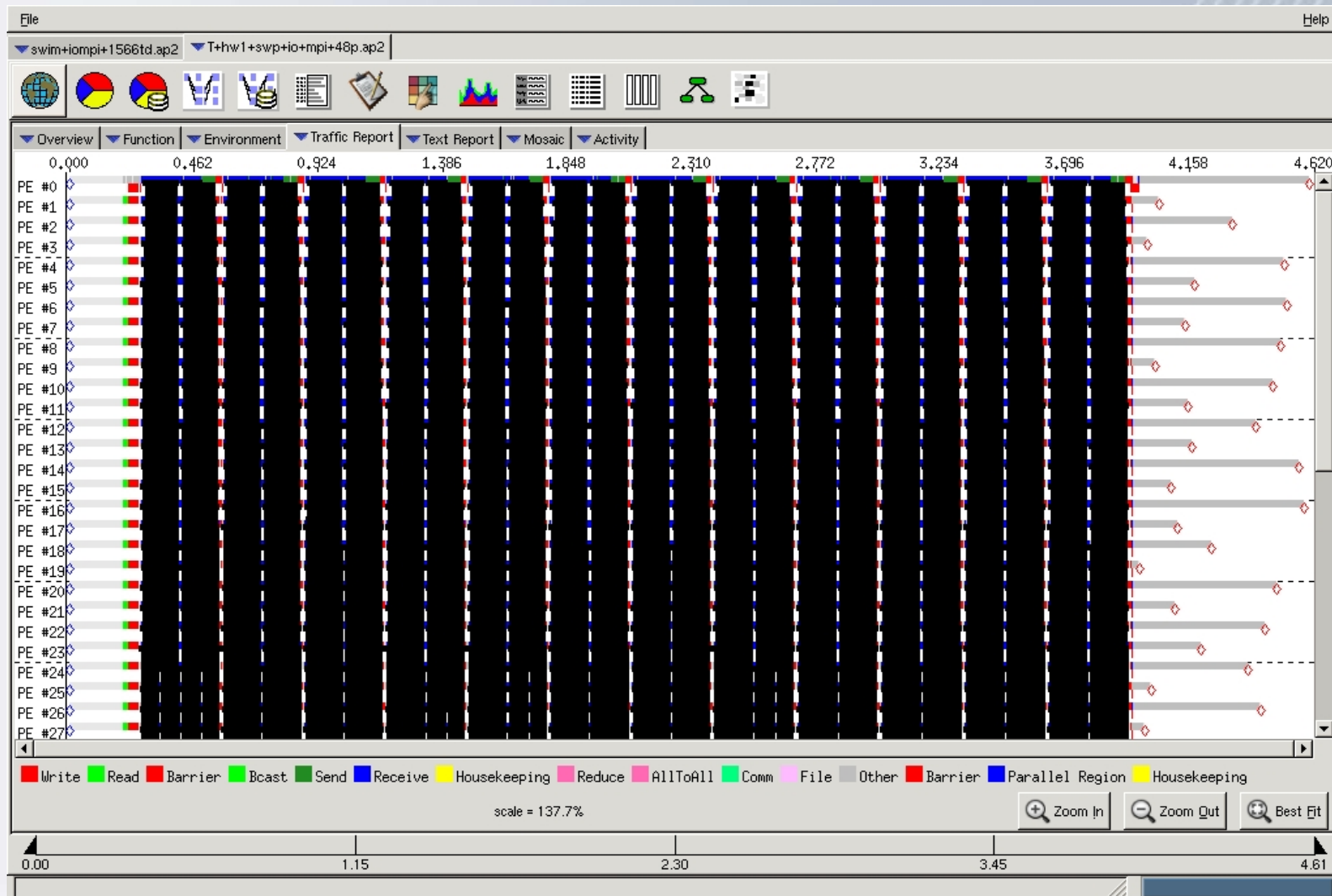
**Environment**

Env Vars System Info Resource Limits Heap Info

PE	Total Used (MB)	Total Free (MB)	Largest Free (MB)	Fragments	
0	97.248817	4065.597900	4065.596680	1614	start
	97.283150	4065.563477	4065.559326	1654	end-start
1	95.571548	4067.275146	4067.274414	1534	start
	95.591171	4067.255615	4067.251953	1592	end-start
2	95.571548	4067.275146	4067.274414	1534	start
	95.591171	4067.255615	4067.251953	1592	end-start
3	95.571548	4067.275146	4067.274414	1534	start
	95.591171	4067.255615	4067.251953	1592	end-start
4	95.571548	4067.275146	4067.274414	1534	start
	95.591171	4067.255615	4067.251953	1592	end-start
5	95.571548	4067.275146	4067.274414	1534	start
	95.591171	4067.255615	4067.251953	1592	end-start
6	95.571548	4067.275146	4067.274414	1534	start
	95.591171	4067.255615	4067.251953	1592	end-start
7	95.571548	4067.275146	4067.274414	1534	start
	95.591171	4067.255615	4067.251953	1592	end-start
8	95.571548	4067.275146	4067.274414	1534	start
	95.591171	4067.255615	4067.251953	1592	end-start
9	95.571548	4067.275146	4067.274414	1534	start
	95.591171	4067.255615	4067.251953	1592	end-start
10	95.571548	4067.275146	4067.274414	1534	start
	95.591171	4067.255615	4067.251953	1592	end-start
11	95.571548	4067.275146	4067.274414	1534	start
	95.591171	4067.255615	4067.251953	1592	end-start

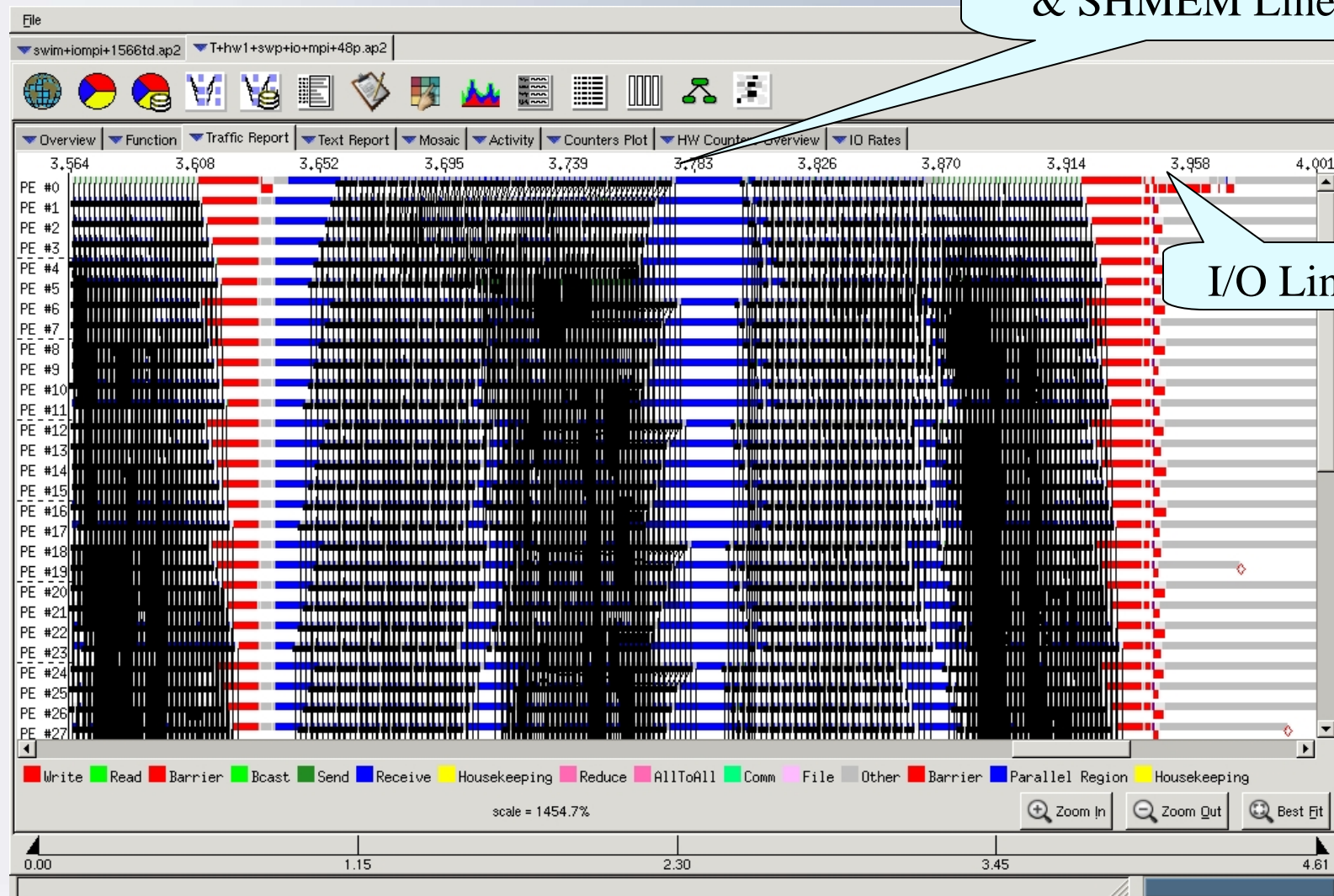
0.00 1.15 2.30 3.45 4.61

# Time Line View (Sweep3D)

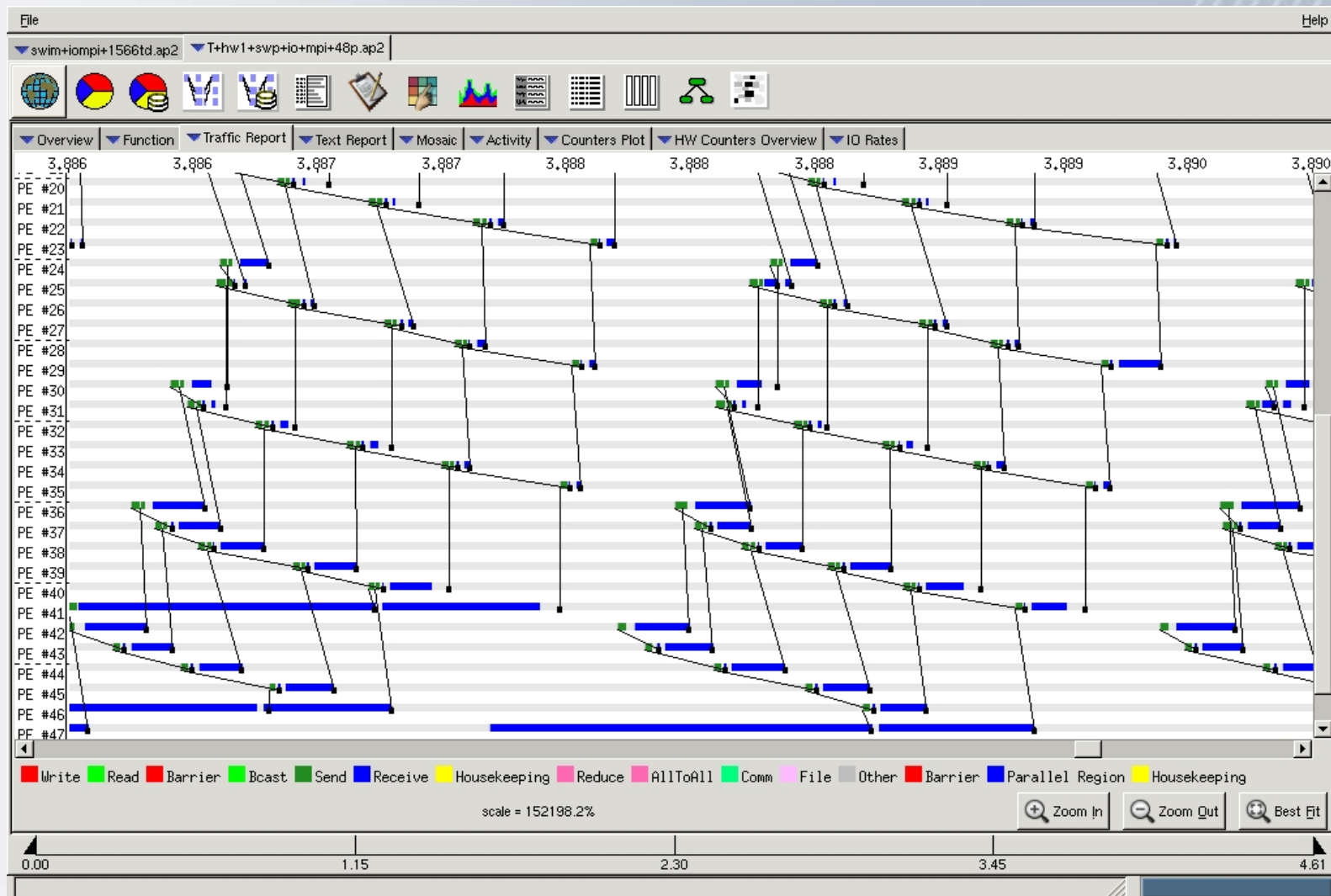




# Time Line View (Zoom)

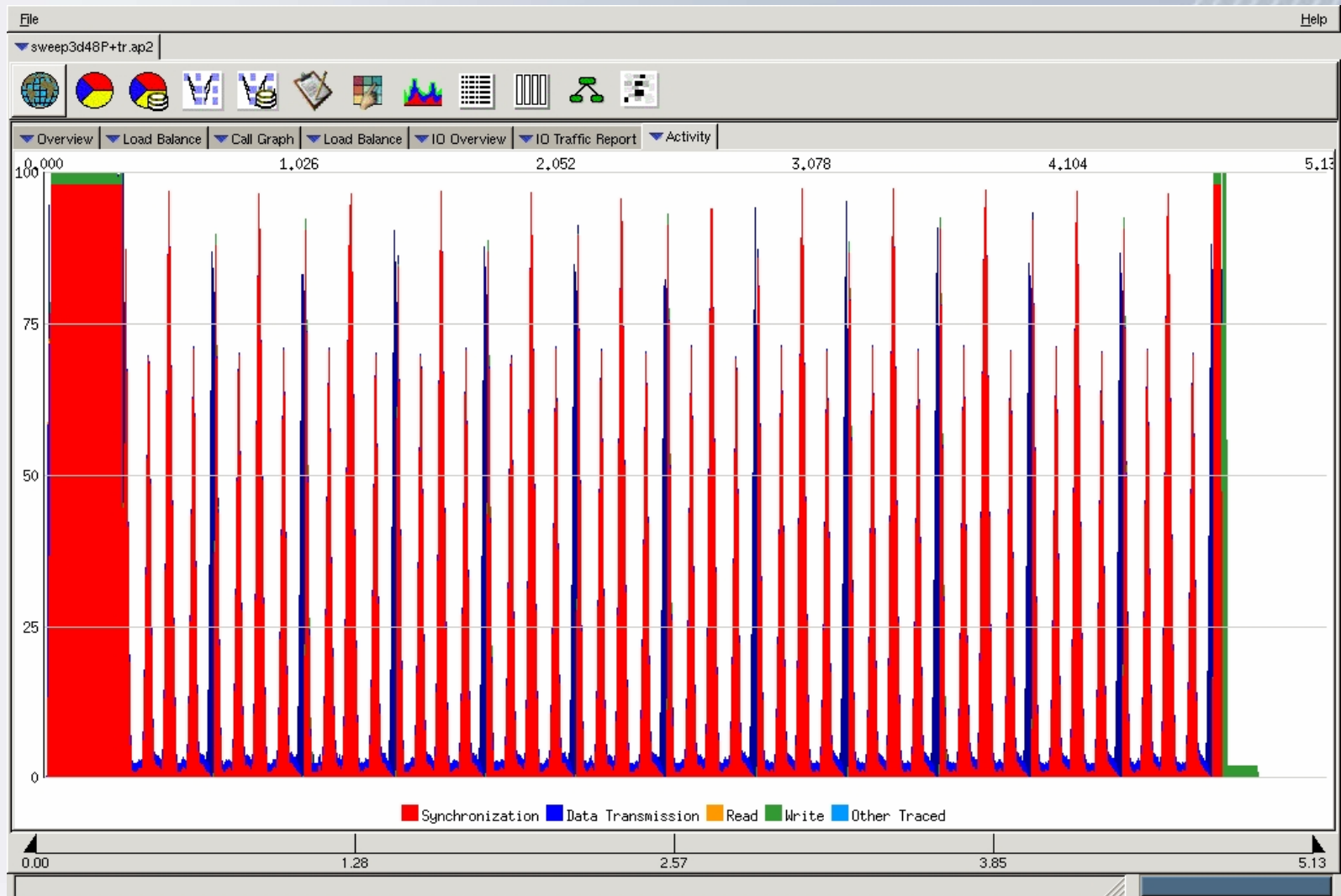


# Time Line View (Fine Grain Zoom)

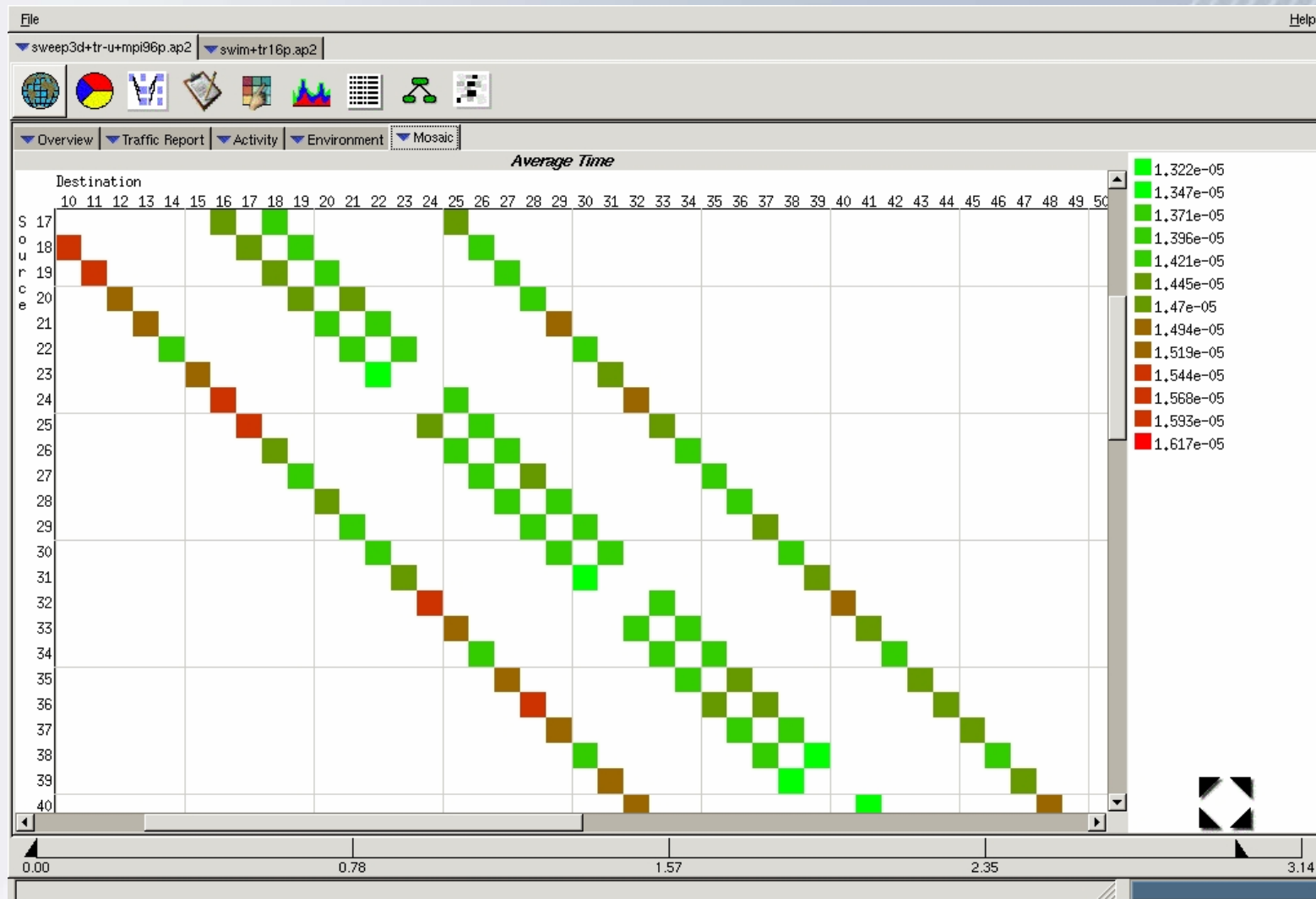




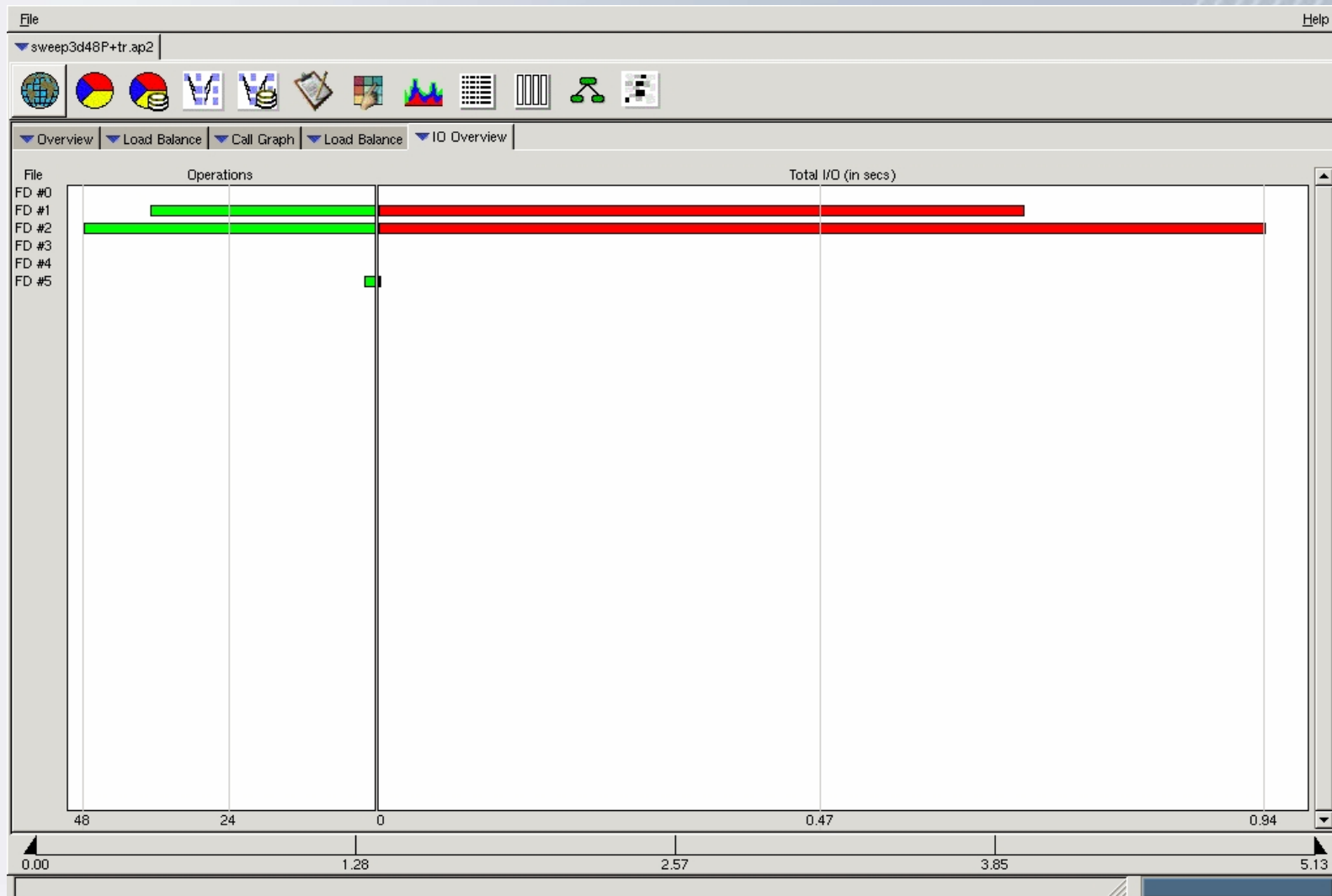
# Activity View



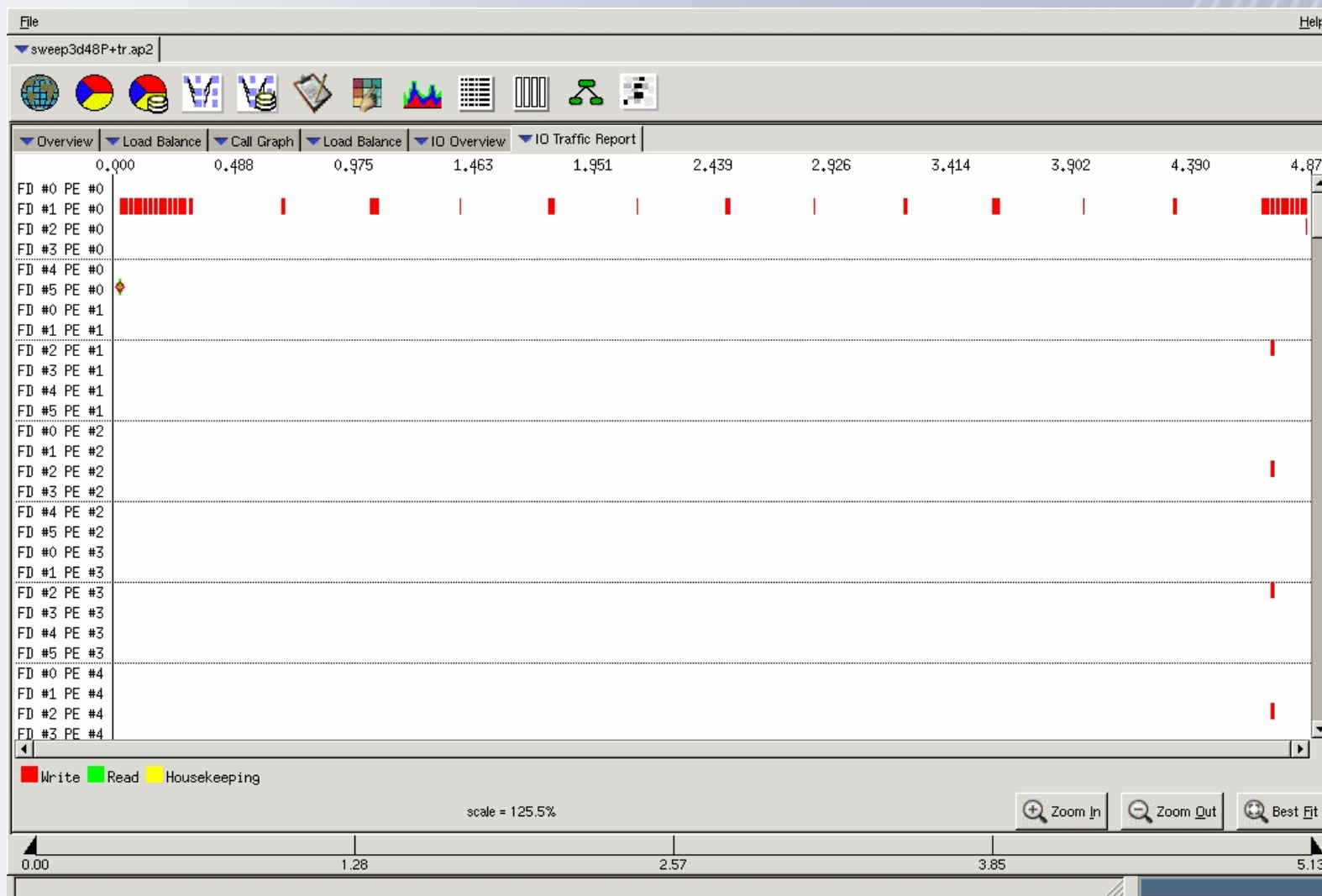
# Pair-wise Communication



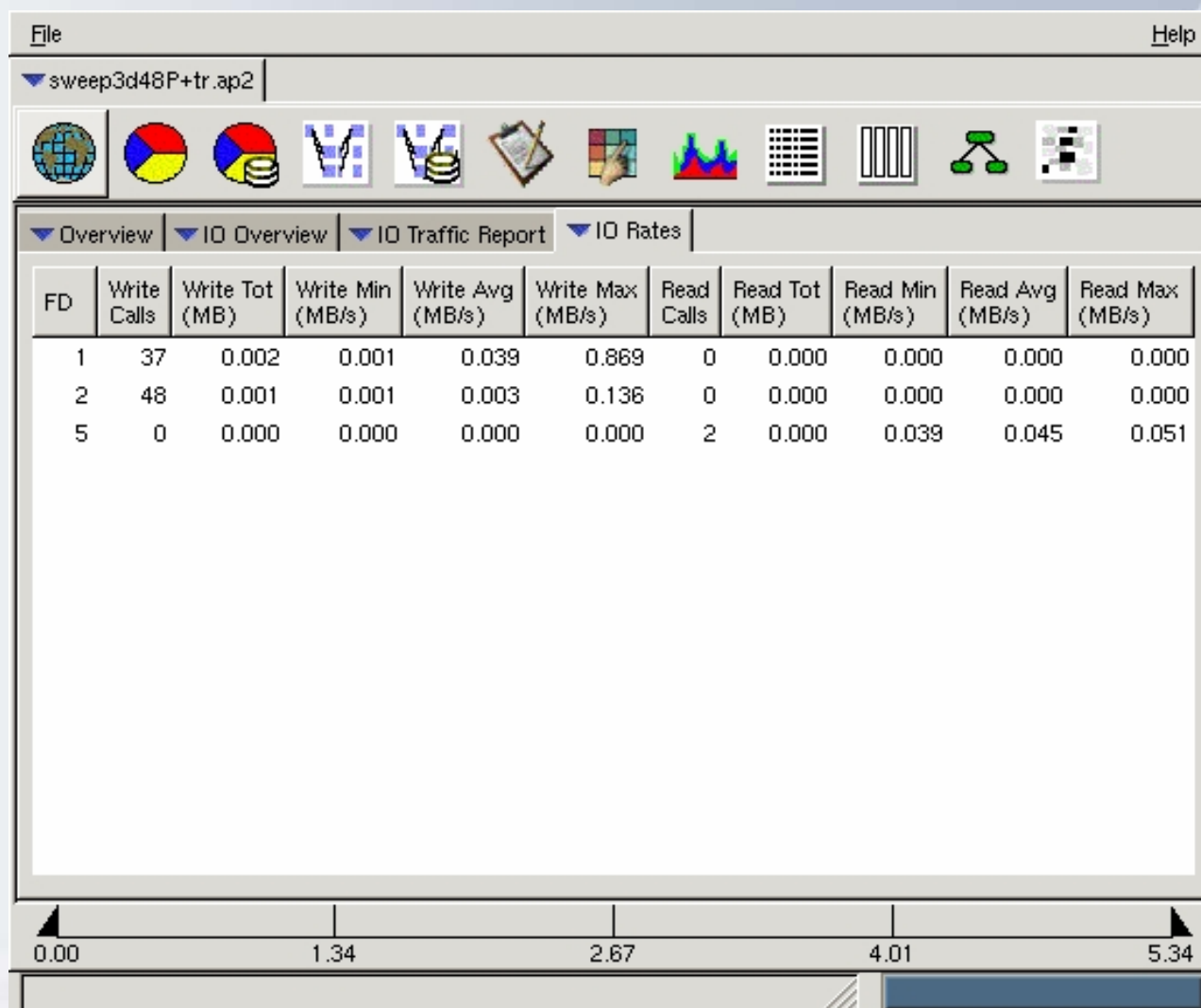
# I/O Overview



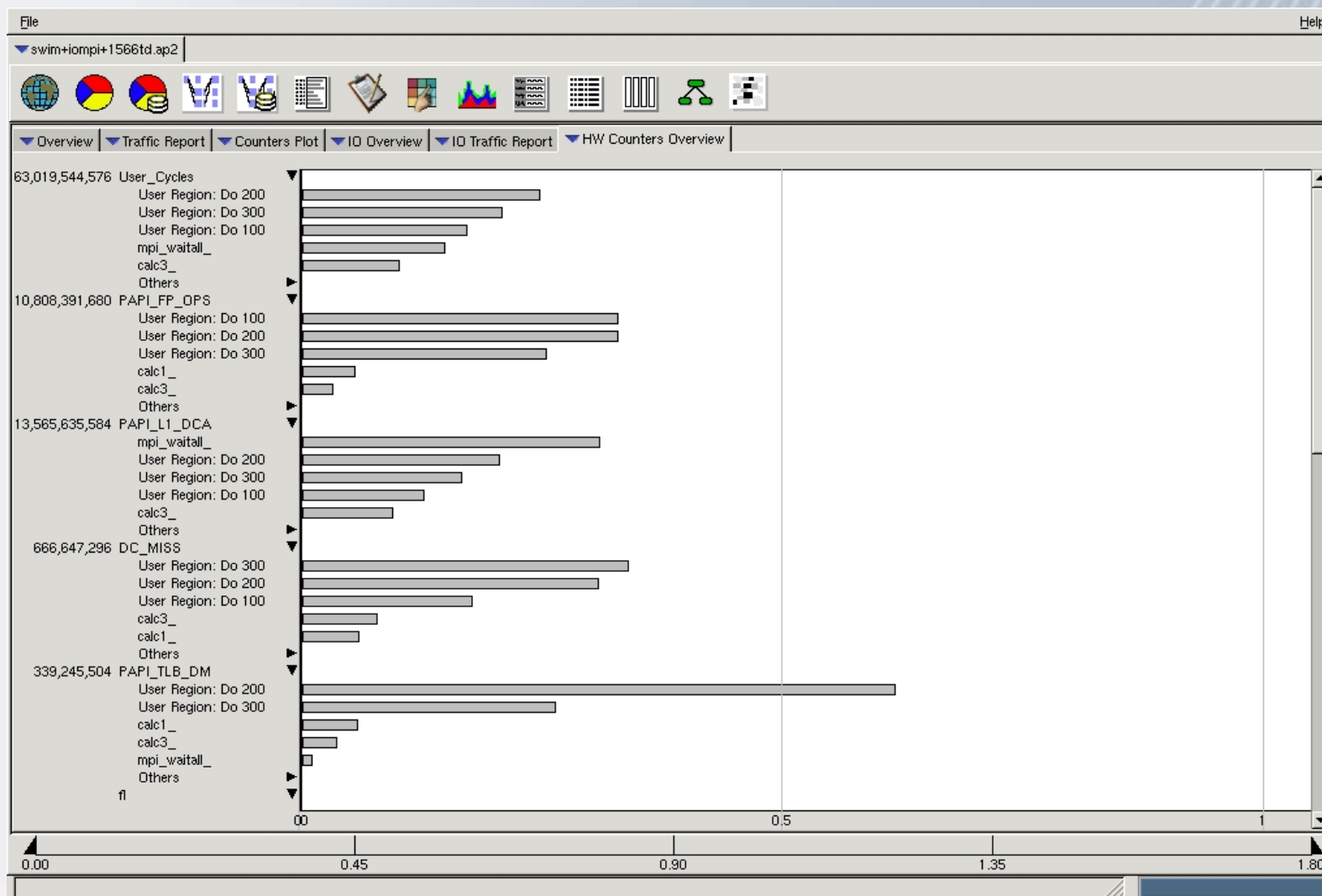
# I/O Traffic Report



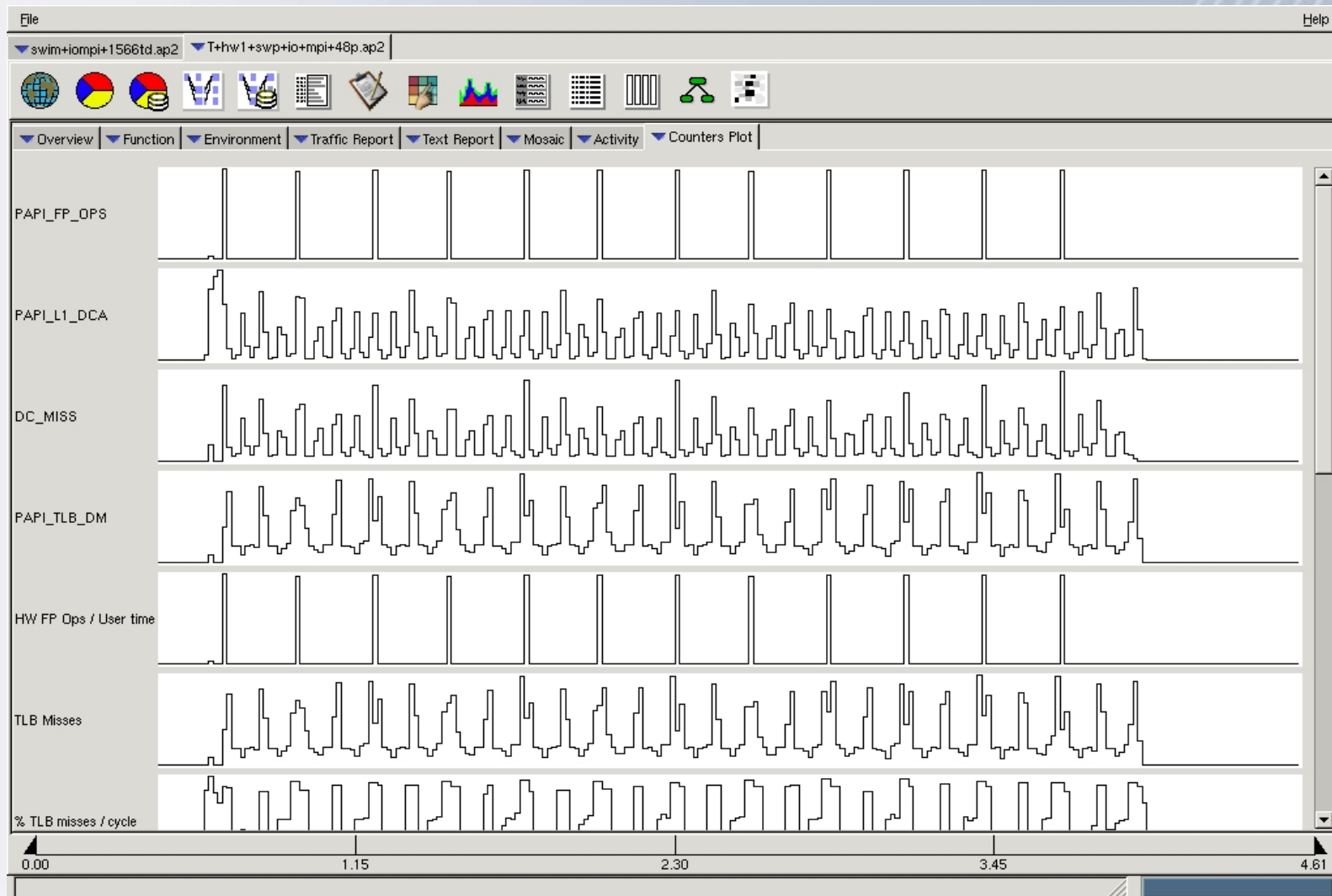
# I/O Rates



# Hardware Counters Overview



# Hardware Counters Time Line



# Controlling Trace File Size

- Several environment variables are available to limit trace files to a reasonable size:
  - **PAT\_RT\_CALLSTACK**
    - Limit the depth to trace the call stack
  - **PAT\_RT\_HWPC**
    - Avoid collecting hardware counters (unset)
  - **PAT\_RT\_RECORD\_PE**
    - Collect trace for a subset of the PEs
  - **PAT\_RT\_TRACE\_FUNCTION\_ARGS**
    - Limit the number of function arguments to be traced
  - **PAT\_RT\_TRACE\_FUNCTION\_LIMITS**
    - Avoid tracing indicated functions
  - **PAT\_RT\_TRACE\_FUNCTION\_MAX**
    - Limit the maximum number of traces generated for all functions for a single process
- Use CrayPat API to toggle trace state (on / off) or to select functions to trace
  - `int PAT_tracing_state (int state)`
  - `int PAT_trace_function (const void *addr, int state)`
- Use the limit built-in command for ksh(1) or csh(1) to control how much disk space the trace file can consume



# Additional API Functions

- int **PAT\_profiling\_state** (int state)
- int **PAT\_record** (int state)
- int **PAT\_sampling\_state** (int state)
- int **PAT\_tracing\_state** (int state)
- int **PAT\_trace\_function** (const void \*addr, int state)
- State can have one of the following:
  - PAT\_STATE\_ON
  - PAT\_STATE\_OFF
  - PAT\_STATE\_QUERY
- int **PAT\_flush\_buffer** (void)

# Parallel Performance Visualization on the Cray XT4

**Questions / Comments**  
**Thank You!**

NERSC  
September 18-20, 2007



Luiz DeRose (ldr@cray.com) © Cray Inc.